

Received June 17, 2020, accepted June 30, 2020, date of publication July 6, 2020, date of current version July 16, 2020.

Digital Object Identifier 10.1109/ACCESS.2020.3007123

Ellipsoidal Subspace Support Vector Data Description

FAHAD SOHRAB¹, (Graduate Student Member, IEEE),

JENNI RAITOHARJU², (Member, IEEE),

ALEXANDROS IOSIFIDIS³, (Senior Member, IEEE),

AND MONCEF GABBOU¹, (Fellow, IEEE)

¹Faculty of Information Technology and Communication Sciences, Tampere University, 33720 Tampere, Finland

²Finnish Environment Institute, 40500 Jyväskylä, Finland

³Department of Engineering, Aarhus University, 8200 Aarhus, Denmark

Corresponding author: Fahad Sohrab (fahad.sohrab@tuni.fi)

This work was supported in part by the National Science Foundation (NSF)-Business Finland Center for Visual and Decision Informatics (CVDI) Project Amalia, and in part by the Business Finland projects VIRPA D and Industrial Data Excellence (INDEX) (Digital, Internet, Materials and Engineering Co-Creation (DIMECC) Industrial Data program).

ABSTRACT In this paper, we propose a novel method for transforming data into a low-dimensional space optimized for one-class classification. The proposed method iteratively transforms data into a new subspace optimized for ellipsoidal encapsulation of target class data. We provide both linear and non-linear formulations for the proposed method. The method takes into account the covariance of the data in the subspace; hence, it yields a more generalized solution as compared to the data description in the subspace by hyperspherical encapsulation of target class data. We propose different regularization terms expressing the class variance in the projected space. We compare the results with classic and recently proposed one-class classification methods and achieve competing results and show clear improvement compared to the other support vector based methods. The proposed method is also noticed to converge much faster than recently proposed Subspace Support Vector Data Description.

INDEX TERMS Anomaly detection, ellipsoidal data description, machine learning, one-class classification, subspace learning.

I. INTRODUCTION

The ability of machines to make a concise description of information requires learning from previous experience. Researchers have been trying to develop techniques for accurately modeling data using supervised and unsupervised learning techniques for many decades. In unsupervised learning techniques, patterns are found without any knowledge of class labels [1]. In supervised learning, labeled training data are used to train models for classifying future instances into different categories [2]. A typical multi-class classification task can be decomposed into several binary classification tasks, where the aim is to decide to which of the two considered classes samples belong to [3]. In binary classification, the data from both classes are used to train a model. One-class classification is conceptually close to binary classification, but the models for classifying future instances are trained using data only from one particular target class [4], [5].

The associate editor coordinating the review of this manuscript and approving it for publication was Shagufta Henna.

In practice, one-class classification is used when data from one of the classes is scarce.

In one-class classification, the class of interest to be modeled is called target or positive class, while samples from the other unknown class(es) are referred to as outliers or negative samples. Numerous attempts have been made to solve one-class classification tasks [6]. The three main approaches for solving one-class classification tasks are density based, reconstruction based, and border based methods [7]. In the density based approach, the description of the target class is based on its density [8], which is usually estimated by using popular density estimation methods such as Parzen density, Gaussian model, or mixture of Gaussians [9]. In reconstruction based approach, some assumptions about the data generating process are made. The underlying function which represents the target class is obtained by fitting a curve over the data by using prior information, such as data clustering characteristics. Self-organizing maps (SOM) [10] and least-squares quantization [11] are classic examples of reconstruction methods. In border based approaches, a model is

created by defining a closed boundary around the target class without estimating its density. One-class Support Vector Machine (OC-SVM) [12] and Support Vector Data Description (SVDD) [13] are among the popular boundary techniques for one-class classification. In OC-SVM, a hyperplane separating the target class is constructed so that the distance of the hyperplane from the origin is maximized. In SVDD, a hypersphere is formed around the target class data by minimizing the volume of hypersphere in a given feature space. Recently, there has been a rising trend to propose approaches based on regression and neural-networks as well [5], [14].

SVDD has been justified over time as a powerful data description method and it has been used in many different application domains for solving one-class classification problems. For example, in [15], SVDD is found to be an excellent choice for solving the problem of identification of freshness of eggs using near infrared spectroscopy (NIR) with an imbalanced number of training samples. In [16], a terrain classification method for ensuring navigation safety of mobile service robots based on SVDD is proposed. To enhance the performance of SVDD, numerous extensions and hybridization techniques have been proposed [8], [17]–[21]. The main extensions of SVDD can be categorized into four main categories. In the first category of extensions, the techniques are focused on manipulating the structure of data, such as associating a confidence coefficient with all training instances which deals with the uncertainty of data [22]. In the second category, the performance is enhanced by proposing new non-linear methods and reducing the complexity of algorithms [23], [24]. Techniques for handling non-stationary data in the context of one-class classification falls in the third category of extensions [25]. In the fourth category, different changes are proposed in the shape of the boundary encapsulating the target data [26].

A popular alternative to the spherical SVDD is Ellipsoidal SVDD (E-SVDD) [26], [27]. E-SVDD forms a unique hyperellipsoid with a minimum volume covering most of the target data. An ellipsoid, unlike a hypersphere, takes into account the difference in variance for each dimension as well as covariance between them. A hypersphere, characterized only by a radius and a center will result in superfluous regions which do not contain any target objects in the input space [28]. Ellipsoids with a minimum volume containing the target data have applications spanning over many different fields. For example, in [29], it is used to detect intrusion in computer networks and, in [30], it is used to estimate the distance between a robot and its surrounding environment for obstacle collision avoidance. An ellipsoid is preferred for heterogeneous data in the input space because its shape is less conservative than a sphere. However, there are some difficulties in kernelizing the algorithms. The kernel trick cannot be applied directly to E-SVDD because its formulation includes outer products rather than inner products [31].

In this paper, we propose a novel subspace learning algorithm for ellipsoidal one-class classification. The proposed method takes into account the covariance of data in the

subspace so that the boundary created around the target class is a better fit. The proposed method finds a projection along with a data description iteratively by minimizing the volume of the hyperellipsoid. We propose different variants of the proposed method by proposing different settings of the regularization term, which takes into account the concentration matrix. We also annexed the regularization term with different settings without taking into account the concentration matrix and report the results. The proposed method is called Ellipsoidal Subspace Support Vector Data Description (ES-SVDD), since it is analogous to Subspace Support Vector Data Description (S-SVDD) [32] but offers more flexibility by using hyperellipsoid instead of hypersphere. Our results show that using hyperellipsoid for data description in the subspace converges faster and produces better results than the data description in a subspace using hypersphere. Further, we see that hyperellipsoid in the subspace optimised for one-class classification provides a better data description as compared to the hyperellipsoid in the original feature space. We also propose a non-linear version of the algorithm by exploiting the non-linear projection trick (NPT) [33].

The rest of the paper is organized as follows. In Section II, we present an overview of related works. In Section III, a detailed derivation of the newly proposed method is presented. In Section IV, we provide and discuss the experimental protocol along with the obtained results and, finally, conclusions are drawn in Section V.

II. BACKGROUND AND RELATED WORK

One-class classification has been studied extensively in recent years and the approaches predominantly focus on data description in a given feature space [7], [13], [22]. On the other hand, feature selection and subspace learning have been an active research area in machine learning, primarily for challenges with data available for all categories [34], [35]. The aim is to avoid the curse of dimensionality in the original feature space by modeling the given data in a lower dimensional space.

In feature selection methods, a subset of representative features is selected by following some criterion [36]–[38]. The two main approaches for feature selection are the *filter* approach and the *wrappers* approach. In the filter approaches, the main focus is on the intrinsic characteristics of the data and they do not take into account any classification algorithm. On the other hand, the wrappers approaches are dependent only on a specific classification algorithm [39].

In subspace learning, the features are transformed from original feature space to a lower-dimensional subspace [40]. Most of the existing subspace learning methods, particularly for anomaly detection, follow three general steps [41], [42]: First, the features are selected randomly by applying random projections to the attributes. Second, classical algorithms are applied locally in each subspace and scores (e.g., voting) are computed. Finally, all the scores are aggregated to compute a global score for classification.

The focus of our paper is to find an optimized subspace for one-class classification. We review the classical one-class classification method, SVDD, in Section II-A and also provide an overview of S-SVDD and graph embedded one-class classifiers in Sections II-B and II-C, respectively.

A. SUPPORT VECTOR DATA DESCRIPTION

Let us denote the data points to be enclosed inside a closed boundary by a matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, $\mathbf{x}_i \in \mathbb{R}^D$, where N is total number of instances and D is dimensionality of data in the original feature space. All the data samples represented by \mathbf{X} belong to the same class.

SVDD finds a spherical boundary around the data by minimizing the volume of a hypersphere enclosing all the target class data:

$$\begin{aligned} \min F(R, \mathbf{a}) &= R^2 \\ \text{s.t. } \|\mathbf{x}_i - \mathbf{a}\|_2^2 &\leq R^2, \quad \forall i \in \{1, \dots, N\}, \end{aligned} \quad (1)$$

where R is the radius of hypersphere and $\mathbf{a} \in \mathbb{R}^D$ is the center of the hypersphere in the given feature space. Slack variables ξ_i , $i = 1, \dots, N$ are introduced for allowing the possibility of data points being outliers, hence the optimization problem changes to

$$\begin{aligned} \min F(R, \mathbf{a}) &= R^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } \|\mathbf{x}_i - \mathbf{a}\|_2^2 &\leq R^2 + \xi_i, \\ \xi_i &\geq 0, \quad \forall i \in \{1, \dots, N\}, \end{aligned} \quad (2)$$

where $C > 0$ is a hyperparameter which controls the trade-off between the volume of the sphere and the amount of data outside the sphere. The Lagrangian dual of (2) reduces to

$$L = \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x}_i - \sum_i \sum_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j, \quad (3)$$

subject to $0 \leq \alpha_i \leq C$. Maximizing (3) gives a set of α_i for corresponding data points. The samples with $\alpha_i > 0$ are the support vectors defining the data description [13]. The samples corresponding to $0 < \alpha_i < C$ lie on the boundary of the hypersphere and those with $\alpha_i = C$ are outliers.

B. SUBSPACE SUPPORT VECTOR DATA DESCRIPTION

In S-SVDD [32], a projection matrix \mathbf{Q} is determined to map data from the original space \mathbb{R}^D to a new optimized lower dimensional space \mathbb{R}^d , $d < D$, so that the data are more suitable for one-class classification:

$$\begin{aligned} \min F(R, \mathbf{a}) &= R^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } \|\mathbf{Q}\mathbf{x}_i - \mathbf{a}\|_2^2 &\leq R^2 + \xi_i, \\ \xi_i &\geq 0, \quad \forall i \in \{1, \dots, N\}, \end{aligned} \quad (4)$$

where $\mathbf{a} \in \mathbb{R}^d$ is the center of the hypersphere in lower d -dimensional space. The method iteratively solves the SVDD in the current subspace to obtain the data description

parameters α_i , $i = 1, \dots, N$, and then updates the subspace projection by optimizing an augmented version of the Lagrangian:

$$L = \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T \mathbf{Q} \mathbf{x}_i - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T \mathbf{Q} \mathbf{x}_j \alpha_j + \beta \psi, \quad (5)$$

where ψ is a regularization term expressing the class variance in the low dimensional space and β is a regularization parameter controlling the importance of the ψ , where

$$\psi = \text{Tr}(\mathbf{Q} \mathbf{X} \boldsymbol{\lambda} \boldsymbol{\lambda}^T \mathbf{X}^T \mathbf{Q}^T), \quad (6)$$

where $\text{Tr}(\cdot)$ is the trace operator and $\boldsymbol{\lambda} \in \mathbb{R}^N$ is a vector controlling the contribution of each training sample. \mathbf{Q} is updated by using the gradient of (5), i.e.,

$$\mathbf{Q} \leftarrow \mathbf{Q} - \eta \Delta L, \quad (7)$$

where η is the learning rate. A non-linear version of S-SVDD employing the kernel trick is also proposed in [32].

C. GRAPH EMBEDDED ONE-CLASS CLASSIFIERS

Graph embedded one-class classifiers constitute extensions of the OC-SVM and SVDD by incorporating generic graph structures in their optimization process. The generic graph structures express geometric data relationships of the target class in the data. For example, Graph Embedded SVDD (GE-SVDD) [17] optimization problem is formulated as

$$\begin{aligned} \min F(R, \mathbf{a}) &= R^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } (\phi(\mathbf{x}_i) - \mathbf{a})^T \mathbf{S}^{-1} (\phi(\mathbf{x}_i) - \mathbf{a}) &\leq R^2 + \xi_i, \\ \xi_i &\geq 0, \quad \forall i \in \{1, \dots, N\}, \end{aligned} \quad (8)$$

where $\phi(\cdot)$ is any non-linear function used for mapping the training samples from the input feature space to the kernel space. The matrix \mathbf{S} contains the geometric data relationships. For example, in PCA, the scatter of training data can be expressed as

$$\mathbf{S} = \frac{1}{N} \Phi \left(\mathbf{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^T \right) \Phi^T = \Phi \mathbf{L} \Phi^T, \quad (9)$$

where $\mathbf{1} \in \mathbb{R}^N$ is a vector containing all values as ones, $\mathbf{I} \in \mathbb{R}^{N \times N}$ is an identity matrix, and Φ is a matrix that contains the training data representations in kernel space.

The Lagrangian of GE-SVDD is

$$\begin{aligned} L &= \sum_{i=1}^N \alpha_i \phi(\mathbf{x}_i)^T \mathbf{S}^{-1} \phi(\mathbf{x}_i) \\ &\quad - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \phi(\mathbf{x}_i)^T \mathbf{S}^{-1} \phi(\mathbf{x}_j) \alpha_j. \end{aligned} \quad (10)$$

It has been shown in [17] that the optimization problem in (10) is equivalent to the problem of SVDD in a transformed feature space.

III. ELLIPSOID SUBSPACE SUPPORT VECTOR DATA DESCRIPTION

Our aim is to find a projection matrix $\mathbf{Q} \in \mathbb{R}^{d \times D}$ to be used for transforming the data to an optimized subspace suitable for one-class classification. In the following analysis, we assume that the data has been centered by setting $\mathbf{X} \leftarrow \mathbf{X} - \boldsymbol{\mu}$, where $\boldsymbol{\mu}$ is the mean of the given training data. The mapping from the original feature space with dimensionality D to a subspace with dimensionality $d \leq D$ is carried out. The mapping is done to transform the data so that it is more suitable to be encapsulated inside an ellipsoid with a minimum volume.

The optimization problem is formulated as

$$\begin{aligned} \min F(R, \mathbf{a}) &= R^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } (\mathbf{Q}\mathbf{x}_i - \mathbf{a})^T \mathbf{E}^{-1} (\mathbf{Q}\mathbf{x}_i - \mathbf{a}) &\leq R^2 + \xi_i, \\ \xi_i &\geq 0, \quad \forall i \in \{1, \dots, N\}, \end{aligned} \quad (11)$$

where \mathbf{a} is the center of the hyperellipsoid and $\mathbf{E} = \mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T$ is the covariance matrix of the data in d -dimensional space. The inverse of covariance matrix \mathbf{E} , also known as the concentration or precision matrix is symmetric and positive definite $\mathbf{E}^{-1} \in \mathbb{R}^{d \times d}$. By defining a new vector $\mathbf{u} = \mathbf{E}^{-\frac{1}{2}}\mathbf{a}$, (11) can be written as

$$\begin{aligned} \min F(R, \mathbf{u}) &= R^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } \|\mathbf{E}^{-\frac{1}{2}}\mathbf{Q}\mathbf{x}_i - \mathbf{u}\|_2^2 &\leq R^2 + \xi_i, \\ \xi_i &\geq 0, \quad \forall i \in \{1, \dots, N\}. \end{aligned} \quad (12)$$

The data in the subspace is represented by

$$\mathbf{y}_i = \mathbf{Q}\mathbf{x}_i, \quad i = 1, \dots, N. \quad (13)$$

The constraints in (12) can be incorporated into its corresponding objective function by using Lagrange multipliers:

$$\begin{aligned} L &= R^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i (R^2 + \xi_i \\ &\quad - (\mathbf{E}^{-\frac{1}{2}}\mathbf{y}_i)^T \mathbf{E}^{-\frac{1}{2}}\mathbf{y}_i + 2\mathbf{u}^T \mathbf{E}^{-\frac{1}{2}}\mathbf{y}_i - \mathbf{u}^T \mathbf{u}) - \sum_{i=1}^N \gamma_i \xi_i \end{aligned} \quad (14)$$

with Lagrange multipliers $\alpha_i \geq 0$ and $\gamma_i \geq 0$.

By setting partial derivatives with respect to R , \mathbf{u} and ξ_i to zero, we get

$$\frac{\partial L}{\partial R} = 0 \Rightarrow \sum_{i=1}^N \alpha_i = 1 \quad (15)$$

$$\frac{\partial L}{\partial \mathbf{u}} = 0 \Rightarrow \mathbf{u} = \sum_{i=1}^N \alpha_i \mathbf{E}^{-\frac{1}{2}}\mathbf{Q}\mathbf{x}_i \quad (16)$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow C - \alpha_i - \xi_i = 0. \quad (17)$$

By substituting (15)-(17) into (14) we get

$$L = \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T \mathbf{E}^{-1} \mathbf{Q}\mathbf{x}_i - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T \mathbf{E}^{-1} \mathbf{Q}\mathbf{x}_j \alpha_j. \quad (18)$$

We can use SVDD to solve (18) for getting α_i values. The concentration matrix \mathbf{E}^{-1} is equivalent to

$$\mathbf{E}^{-1} = (\mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T)^{-1}. \quad (19)$$

By putting (19) in (18) we get

$$\begin{aligned} L &= \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T (\mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T)^{-1} \mathbf{Q}\mathbf{x}_i \\ &\quad - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T (\mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T)^{-1} \mathbf{Q}\mathbf{x}_j \alpha_j. \end{aligned} \quad (20)$$

We add an extra term Υ to (20) as a regularization term expressing the class variance in the projected space, also taking into account the concentration matrix. Hence, (20) now becomes

$$\begin{aligned} L &= \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T (\mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T)^{-1} \mathbf{Q}\mathbf{x}_i \\ &\quad - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \mathbf{x}_i^T \mathbf{Q}^T (\mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T)^{-1} \mathbf{Q}\mathbf{x}_j \alpha_j + \beta \Upsilon, \end{aligned} \quad (21)$$

where β controls the importance of regularization term and is used as a hyperparameter. Υ is defined as follows:

$$\Upsilon = \text{Tr}(\mathbf{E}^{-\frac{1}{2}}\mathbf{Q}\mathbf{X}\boldsymbol{\lambda}\boldsymbol{\lambda}^T\mathbf{X}^T\mathbf{Q}^T\mathbf{E}^{-\frac{1}{2}}), \quad (22)$$

where $\boldsymbol{\lambda}$ can take three different forms. In the first form, all elements in $\boldsymbol{\lambda}$ take the value of 1 and, hence, all the samples are used to describe the covariance of the class. In the second form, $\boldsymbol{\lambda}$ is replaced by $\boldsymbol{\alpha}$, which means that the samples belonging to the boundary and outside the boundary are used to describe the covariance of the class. In the third form, the λ_i values are replaced by α_i values of the samples belonging to the boundary and zero for other instances. The first, second and third forms of the regularization terms are expressed as Υ_1 , Υ_2 , and Υ_3 hereinafter.

In our experiments, we also consider the regularization term expressing the class variance in the projected space without taking into account the concentration matrix. This is achieved by replacing the covariance matrix \mathbf{E} with the identity matrix \mathbf{I} in (22). By doing so, the regularization term Υ becomes equivalent to ψ as described in (6). Analogous to regularization term Υ , ψ can also take different forms by changing $\boldsymbol{\lambda}$ and similarly hereinafter we refer to all those cases by ψ_1 , ψ_2 and ψ_3 . The methods used with ψ and Υ are denoted by ES-SVDD ψ_m and ES-SVDD Υ_m ($m = 1, 2, 3$), respectively. We refer to the case, where no regularization term is used in ES-SVDD, as ES-SVDD $\psi_0\Upsilon_0$.

Equation (21) can be further simplified and written as

$$L = \text{Tr}((\mathbf{Q}\mathbf{X}\mathbf{X}^T\mathbf{Q}^T)^{-1}\mathbf{Q}\mathbf{X}(\mathbf{A} - \boldsymbol{\alpha}\boldsymbol{\alpha}^T)\mathbf{X}^T\mathbf{Q}^T) + \beta \Upsilon, \quad (23)$$

where \mathbb{A} is a diagonal matrix having α_i values in its diagonal and $\boldsymbol{\alpha}$ is a vector of α_i 's. We use gradient of (23) to update the projection matrix. The gradient can be solved using identity 126 in [43]:

$$\Delta L = 2\mathbf{E}^{-1}\mathbf{Q}\mathbf{X}(\mathbb{A} - \boldsymbol{\alpha}\boldsymbol{\alpha}^\top)\mathbf{X}^\top - 2\mathbf{E}^{-1}\mathbf{Q}\mathbf{X}(\mathbb{A} - \boldsymbol{\alpha}\boldsymbol{\alpha}^\top)\mathbf{X}^\top\mathbf{Q}^\top\mathbf{E}^{-1}\mathbf{Q}\mathbf{X}\mathbf{X}^\top + \beta\Delta\Upsilon, \quad (24)$$

where

$$\Delta\Upsilon = 2\mathbf{E}^{-1}\mathbf{Q}\mathbf{X}\boldsymbol{\lambda}\boldsymbol{\lambda}^\top\mathbf{X}^\top - 2\mathbf{E}^{-1}\mathbf{Q}\mathbf{X}\boldsymbol{\lambda}\boldsymbol{\lambda}^\top\mathbf{X}^\top\mathbf{Q}^\top\mathbf{E}^{-1}\mathbf{Q}\mathbf{X}\mathbf{X}^\top. \quad (25)$$

When ψ is used as a regularization term, we use $\Delta\psi$ instead of $\Delta\Upsilon$ in (24):

$$\Delta\psi = 2\mathbf{Q}\mathbf{X}\boldsymbol{\lambda}\boldsymbol{\lambda}^\top\mathbf{X}^\top. \quad (26)$$

We obtain an optimised data projection matrix along with optimised data description in a two-step iterative process. In the first step, the α_i values are computed by maximizing (18). In the second step, \mathbf{Q} is updated through the gradient descent after computing the gradient by using (23). In order to obtain an orthogonal projection, we impose the orthogonality constraint $\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}$. We orthogonalize and normalize \mathbf{Q} during the two-step iterative process. Algorithm 1 presents the whole algorithm.

Algorithm 1 Linear ES-SVDD Optimization

Input : \mathbf{X} , β , η , d , C , k_{max}

Output: \mathbf{Q} , R , $\boldsymbol{\alpha}$

Random initialization of \mathbf{Q} ;

Initialize $k = 1$;

while $k < k_{max}$ **do**

 Compute concentration matrix \mathbf{E}^{-1} using (19);

 Solve α_i , $i = 1, \dots, N$ with SVDD using (18);

 Calculate ΔL using (24);

 Update $\mathbf{Q} \leftarrow \mathbf{Q} - \eta\Delta L$;

 Orthogonalize \mathbf{Q} using QR decomposition;

 Row normalize \mathbf{Q} using l_2 norm;

$k \leftarrow k + 1$

end

// Data description in the optimized subspace

Compute concentration matrix \mathbf{E}^{-1} using (19)

Calculate α_i , $i = 1, \dots, N$ with SVDD using (18);

A. NON-LINEAR ELLIPSOIDAL SUBSPACE SUPPORT VECTOR DATA DESCRIPTION

The non-linear ellipsoidal subspace SVDD is not trivial, because the kernel trick cannot be applied directly due to the outer products involved in its derivation. To avoid this problem, we follow the NPT based solution described below [33]. We first compute a noncentered kernel matrix $\mathbf{K} = \Phi^\top\Phi$ using the radial basis function kernel as

$$\mathbf{K}_{ij} = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right), \quad (27)$$

where σ is a hyperparameter scaling the distance between \mathbf{x}_i and \mathbf{x}_j . The kernel matrix is centered as

$$\hat{\mathbf{K}} = (\mathbf{I} - \mathbf{J})\mathbf{K}(\mathbf{I} - \mathbf{J}), \quad (28)$$

where $\mathbf{J} \in \mathbb{R}^{N \times N}$ is a matrix defined as

$$\mathbf{J} = \frac{1}{N}\mathbf{1}\mathbf{1}^\top. \quad (29)$$

The centered kernel matrix is decomposed by using eigendecomposition:

$$\hat{\mathbf{K}} = \mathbf{U}\mathbf{A}\mathbf{U}^\top, \quad (30)$$

where \mathbf{A} contains the non-negative eigenvalues of the centered kernel matrix in its diagonal and the columns of \mathbf{U} contain the corresponding eigenvectors. Finally, the data in the reduced dimensional kernel space is obtained as

$$\Phi = (\mathbf{A}^{\frac{1}{2}})^+ \mathbf{U}^+ \hat{\mathbf{K}}, \quad (31)$$

where $+$ sign in the superscript denotes the pseudo-inverse.

After applying NPT, we continue by considering Φ as our input data. This allows us to use the linear E-SVDD formulation to obtain a non-linear transformation.

B. TEST PHASE

During the test phase of the linear case, a test instance \mathbf{x}_* is first mapped to the optimized lower d -dimensional space as

$$\mathbf{y}_* = \mathbf{Q}\mathbf{x}_*. \quad (32)$$

The decision to classify the instance as target or outlier is taken on the basis of its distance from the center of data description in the d -dimensional space. The distance is calculated as follows:

$$\|\mathbf{E}^{-\frac{1}{2}}\mathbf{y}_* - \mathbf{u}\|_2^2 = (\mathbf{E}^{-\frac{1}{2}}\mathbf{y}_*)^\top \mathbf{E}^{-\frac{1}{2}}\mathbf{y}_* - 2(\mathbf{E}^{-\frac{1}{2}}\mathbf{y}_*)^\top \mathbf{u} + \mathbf{u}^\top \mathbf{u}, \quad (33)$$

where \mathbf{u} can be solved with (16). If $\|\mathbf{E}^{-\frac{1}{2}}\mathbf{y}_* - \mathbf{u}\|_2^2 \leq R^2$, the test instance is classified as positive, as it will fall inside the boundary of the data description. The test instance is classified as negative if $\|\mathbf{E}^{-\frac{1}{2}}\mathbf{y}_* - \mathbf{u}\|_2^2 > R^2$. The threshold R^2 for taking the decision is calculated as follows:

$$R^2 = (\mathbf{E}^{-\frac{1}{2}}\mathbf{s})^\top \mathbf{E}^{-\frac{1}{2}}\mathbf{s} - 2\mathbf{u}^\top \mathbf{s} + \mathbf{u}^\top \mathbf{u}, \quad (34)$$

where \mathbf{s} is any support vector with $0 < \alpha_i < C$.

During the test phase for non-linear ES-SVDD, we use NPT by first computing the kernel vector as

$$\mathbf{k}_* = \Phi^\top \phi(\mathbf{x}_*). \quad (35)$$

The kernel vector is then centered as

$$\hat{\mathbf{k}}_* = (\mathbf{I} - \mathbf{J})[\mathbf{k}_* - \frac{1}{N}\mathbf{K}\mathbf{1}]. \quad (36)$$

The centered kernel vector is then mapped to

$$\phi_* = (\Phi^\top)^+ \hat{\mathbf{k}}_* \quad (37)$$

We now consider ϕ_* as the test input \mathbf{x}_* and follow all the steps described for the linear test.

TABLE 1. Datasets used in the experiments.

Abbreviation	Dataset Name (Target Class)	Total Samples	Target Samples	D
S-K	Seeds (Kama)	210	70	7
S-R	Seeds (Rosa)	210	70	7
S-C	Seeds (Canadian)	210	70	7
QB-B	Qualitative bankruptcy (bankruptcy)	250	107	6
QB-N	Qualitative bankruptcy (non-bankruptcy)	250	143	6
SH-H	Somerville happiness (happy)	143	77	6
SH-U	Somerville happiness (un-happy)	143	66	6
I-S	Iris (Setosa)	150	50	4
I-VC	Iris (Versicolor)	150	50	4
I-V	Iris (Virginica)	150	50	4
IS-B	Ionosphere (bad)	351	126	34
IS-G	Ionosphere (good)	351	225	34
SR-R	Sonar (rock)	208	97	60
SR-M	Sonar (mines)	208	111	60

IV. EXPERIMENTS

A. DATASETS AND EXPERIMENTAL SETUP

We evaluated the proposed and competing methods over different datasets downloaded from UCI machine learning repository [44]. Since one-class classification methods inherently are suited for binary (target and outliers) imbalanced classification problems, we converted the datasets to one-class datasets by considering a single class in a dataset at a time as the target class and all other classes as outliers. Naturally, only the target class samples were used for training the models, while all the classes were considered in the validation and test phases. The total number of samples, number of target class samples, and number of dimensions in the original feature space are given in Table 1.

In each dataset, 70% of the data was used for training and the remaining 30% for testing. The train and test sets were selected randomly by keeping the proportions of classes similar to the full dataset. Each experiment was repeated five times using different random train/test splits, while the same five splittings were used for all the considered methods. We report the average test performance over the five splittings. During training, a 5-fold cross-validation technique was used to select the best hyperparameters with the best evaluation score. We used only the training sets for selecting the hyperparameters. We used Geometric mean ($Gmean$) as the evaluation metric for all the methods. $Gmean$ is defined as

$$Gmean = \sqrt{tpr \times tnr}, \quad (38)$$

where tpr is true positive rate (also known as sensitivity) and tnr is true negative rate (also known as specificity). For the proposed ES-SVDD method, we chose the hyperparameters from the following values

- $\beta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4\}$,
- $C \in \{0.01, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6\}$,
- $\sigma \in \{10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3\}$,
- $d \in \{1, 2, 3, 4, 5, 10, 20, 50, 100\}$,
- $\eta \in \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$.

For all the competing methods, the hyperparameters corresponding to ES-SVDD hyperparameters were selected from

the above values. For other hyperparameters, the same ranges were used as provided in the corresponding work or stated otherwise. We used the target class samples of the full training set with the optimal hyperparameters for the final training and then tested with the test set.

We compared the proposed ES-SVDD with other support vector (SV)-based and non-SV-based methods. The SV-based one-class classification methods essentially create a model by defining a boundary. The SV-based methods used for comparison were S-SVDD [32], OC-SVM [12], SVDD [13], and E-SVDD. The non-SV-based methods used for comparison were density-based, reconstruction-based, and regression-based one-class classification approaches. The density-based methods used for comparison were Parzen density-based data description [7] and Gaussian density-based data description [7]. As reconstruction-based methods, we used SOM data description [7] and K-means data description [7]. The regression-based method used for comparison was Graph Embedded One-Class Extreme Learning Machines (GE-OC-ELM) which exploits geometric class information [5].

We used maximum likelihood estimation for finding the optimum smoothing parameter in the Parzen density-based data description. The grid-size in SOM was fixed to $5\sqrt{N_t}$, where N_t is the size of training data for a given dataset [45]. We chose the number of clusters (N_c) for K-means from $N_c = \{1, 2, 3\}$ and report the best outcome. For non-linear methods, we employed NPT for ES-SVDD and S-SVDD, kernel whitening for Gaussian data description [46], and the kernel trick for other methods. Since the closest counterpart of the proposed method is S-SVDD and different regularization terms for S-SVDD were proposed [32], we report the results with all the previously proposed variants of S-SVDD. We used LIBSVM [47] toolbox implementation for OC-SVM and SVDD and DD-toolbox [48] for SOM, K-means data description, Parzen density, and Gaussian density-based data description. The implementation of GE-OC-ELM is publicly available.¹ The proposed ES-SVDD, along with S-SVDD

¹<https://sites.google.com/view/iosifidis/codes-and-datasets>

TABLE 2. Gmean results for linear methods over different datasets.

Dataset	S-K	S-R	S-C	Av.	QB-B	QB-N	Av.	SH-H	SH-U	Av.
ES-SVDD Υ_0	0.83	0.91	0.77	0.84	0.84	0.26	0.55	0.41	0.51	0.46
ES-SVDD ψ_1	0.83	0.89	0.89	0.87	0.76	0.46	0.61	0.47	0.47	0.47
ES-SVDD ψ_2	0.82	0.89	0.87	0.86	0.85	0.18	0.51	0.53	0.51	0.52
ES-SVDD ψ_3	0.79	0.90	0.87	0.86	0.90	0.23	0.57	0.46	0.35	0.40
ES-SVDD Υ_1	0.82	0.92	0.90	0.88	0.85	0.32	0.58	0.46	0.52	0.49
ES-SVDD Υ_2	0.84	0.91	0.91	0.89	0.81	0.30	0.56	0.55	0.53	0.54
ES-SVDD Υ_3	0.85	0.88	0.88	0.87	0.87	0.33	0.60	0.49	0.47	0.48
S-SVDD ψ_0	0.79	0.86	0.81	0.82	0.72	0.50	0.61	0.49	0.48	0.48
S-SVDD ψ_1	0.72	0.76	0.77	0.75	0.85	0.34	0.59	0.46	0.46	0.46
S-SVDD ψ_2	0.81	0.82	0.77	0.80	0.75	0.40	0.58	0.47	0.48	0.48
S-SVDD ψ_3	0.80	0.93	0.75	0.83	0.72	0.41	0.57	0.49	0.46	0.48
SVDD	0.82	0.92	0.86	0.87	0.83	0.04	0.43	0.54	0.48	0.51
E-SVDD	0.80	0.87	0.86	0.84	0.96	0.20	0.58	0.54	0.41	0.48
OC-SVM	0.43	0.46	0.58	0.49	0.46	0.55	0.51	0.45	0.42	0.44
Non-support-vector-based methods										
K-means	0.86	0.94	0.91	0.90	0.71	0.41	0.56	0.56	0.39	0.47
Parzen	0.49	0.33	0.58	0.47	0.98	0.60	0.79	0.58	0.43	0.50

Dataset	I-S	I-VC	I-V	Av.	IS-B	IS-G	Av.	SR-R	SR-M	Av.
ES-SVDD ψ_0	0.64	0.75	0.70	0.70	0.16	0.89	0.52	0.50	0.64	0.57
ES-SVDD ψ_1	0.92	0.86	0.77	0.85	0.52	0.85	0.69	0.50	0.56	0.53
ES-SVDD ψ_2	0.87	0.82	0.79	0.83	0.31	0.87	0.59	0.48	0.67	0.58
ES-SVDD ψ_3	0.93	0.87	0.71	0.84	0.35	0.89	0.62	0.48	0.65	0.57
ES-SVDD Υ_1	0.85	0.84	0.86	0.85	0.26	0.87	0.57	0.47	0.67	0.57
ES-SVDD Υ_2	0.96	0.83	0.74	0.84	0.31	0.89	0.60	0.47	0.65	0.56
ES-SVDD Υ_3	0.80	0.85	0.79	0.81	0.35	0.90	0.62	0.49	0.65	0.57
S-SVDD ψ_0	0.87	0.75	0.64	0.76	0.16	0.75	0.46	0.37	0.37	0.37
S-SVDD ψ_1	0.88	0.81	0.75	0.81	0.50	0.71	0.61	0.44	0.36	0.40
S-SVDD ψ_2	0.87	0.84	0.58	0.76	0.43	0.72	0.58	0.46	0.40	0.43
S-SVDD ψ_3	0.81	0.68	0.63	0.71	0.27	0.66	0.46	0.46	0.41	0.43
SVDD	0.94	0.90	0.89	0.91	0.04	0.73	0.39	0.50	0.52	0.51
E-SVDD	0.89	0.88	0.86	0.88	0.33	0.00	0.17	0.00	0.00	0.00
OC-SVM	0.50	0.52	0.39	0.47	0.47	0.45	0.46	0.44	0.52	0.48
Non-support-vector-based methods										
K-means	0.94	0.92	0.89	0.91	0.37	0.88	0.63	0.49	0.68	0.58
Parzen	0.85	0.68	0.79	0.77	0.32	0.25	0.28	0.00	0.00	0.00

and E-SVDD, was implemented by the authors using Matlab by leveraging LIBSVM.

B. EXPERIMENTAL RESULTS AND DISCUSSION

In Tables 2 and 3, we report the average test results for each dataset for the linear and non-linear cases, respectively. In each experiment, a single class was selected as the target class and the rest of the data as outliers (see Table 1). We also report the average performance of the proposed and competing methods in the average (Av.) column by averaging the results for a given dataset. For example, the performance over S-K, S-R, and S-C is averaged and provided in the Av. column as the overall performance for Seeds dataset. In this way, we can get an idea of the overall performance for each algorithm over the full dataset. For ES-SVDD and S-SVDD, we report the test results after 10 training iterations.

When compared to SV-based methods, our proposed methods achieved the best average results on all but Iris dataset in the linear case and on half of the datasets in the non-linear case. We note that the average results for the non-linear methods are generally better than those of the linear ones for the majority of the datasets. Overall, the proposed (linear and non-linear) methods achieved the best average results in 4 out of 6 datasets among the SV-based methods. In general,

the best performing methods vary for different datasets, but we can see that there is no case, where the proposed method would fail completely, unlike most of the competing methods. In the linear case, other SV-based competing methods outperformed ES-SVDD only with Iris dataset, which has the lowest original dimensionality and also a low number of samples. Also in the non-linear case, other SV-based methods outperformed ES-SVDD most clearly on the 2 smallest datasets. Thus, it seems that the proposed method is more beneficial when the data dimensionality is higher.

When compared with also non-SV-based methods, we see that the proposed method gave the best average performance on 3 out of 6 datasets in the linear case. In the non-linear case, the ranking of the methods varies more and only GE-SVM achieved the best average results on more than one (2) datasets. The proposed method outperformed the other methods on Ionosphere dataset, which is the largest considered dataset. Furthermore, the stable performance of the proposed method makes it a viable solution also in the non-linear case.

Comparing regularization terms for linear ES-SVDD, we notice that ES-SVDD performed better in majority of cases with regularization term Υ_2 which uses samples belonging to the boundary and outside the boundary to describe the

TABLE 3. Gmean results for non-linear methods over different datasets.

Dataset	S-K	S-R	S-C	Av.	QB-B	QB-N	Av.	SH-H	SH-U	Av.
ES-SVDD $\psi_0 \Upsilon_0$	0.78	0.88	0.93	0.87	0.83	0.61	0.72	0.52	0.42	0.47
ES-SVDD ψ_1	0.80	0.88	0.88	0.85	0.80	0.34	0.57	0.51	0.42	0.47
ES-SVDD ψ_2	0.80	0.90	0.93	0.87	0.90	0.35	0.62	0.52	0.33	0.42
ES-SVDD ψ_3	0.82	0.86	0.72	0.80	0.89	0.64	0.76	0.52	0.45	0.49
ES-SVDD Υ_1	0.85	0.92	0.91	0.89	0.87	0.47	0.67	0.47	0.38	0.43
ES-SVDD Υ_2	0.82	0.88	0.89	0.86	0.84	0.68	0.76	0.52	0.34	0.43
ES-SVDD Υ_3	0.85	0.88	0.91	0.88	0.87	0.54	0.71	0.52	0.45	0.49
S-SVDD ψ_0	0.74	0.74	0.83	0.77	0.23	0.49	0.36	0.45	0.29	0.37
S-SVDD ψ_1	0.71	0.78	0.81	0.77	0.11	0.08	0.09	0.39	0.32	0.35
S-SVDD ψ_2	0.72	0.85	0.83	0.80	0.36	0.37	0.37	0.47	0.32	0.40
S-SVDD ψ_3	0.60	0.76	0.76	0.71	0.36	0.40	0.38	0.46	0.29	0.37
SVDD	0.86	0.91	0.88	0.88	0.88	0.55	0.71	0.54	0.48	0.51
E-SVDD	0.84	0.85	0.85	0.85	0.96	0.51	0.74	0.55	0.42	0.48
GE-SVDD	0.84	0.90	0.76	0.83	0.94	0.17	0.55	0.54	0.47	0.50
OC-SVM	0.79	0.60	0.63	0.67	0.67	0.52	0.59	0.57	0.48	0.52
GE-SVM	0.83	0.88	0.89	0.87	0.88	0.58	0.73	0.57	0.42	0.50
Non-support-vector-based methods										
SOM	0.80	0.90	0.89	0.86	0.79	0.37	0.58	0.28	0.26	0.27
Gaussian	0.85	0.95	0.94	0.91	0.63	0.46	0.54	0.52	0.42	0.47
GE-OC-ELM	0.85	0.93	0.89	0.89	1.00	0.80	0.90	0.31	0.31	0.31
Dataset	I-S	I-VC	I-V	Av.	IS-B	IS-G	Av.	SR-R	SR-M	Av.
ES-SVDD $\psi_0 \Upsilon_0$	0.93	0.84	0.86	0.88	0.44	0.89	0.67	0.41	0.67	0.54
ES-SVDD ψ_1	0.94	0.81	0.74	0.83	0.71	0.90	0.80	0.48	0.55	0.51
ES-SVDD ψ_2	0.91	0.87	0.83	0.87	0.31	0.87	0.59	0.47	0.66	0.56
ES-SVDD ψ_3	0.89	0.84	0.74	0.82	0.32	0.88	0.60	0.47	0.66	0.57
ES-SVDD Υ_1	0.81	0.89	0.70	0.80	0.47	0.86	0.67	0.53	0.65	0.59
ES-SVDD Υ_2	0.91	0.83	0.81	0.85	0.68	0.86	0.77	0.47	0.70	0.58
ES-SVDD Υ_3	0.94	0.88	0.83	0.88	0.45	0.85	0.65	0.41	0.70	0.55
S-SVDD ψ_0	0.92	0.85	0.78	0.85	0.24	0.53	0.38	0.43	0.41	0.42
S-SVDD ψ_1	0.89	0.89	0.63	0.80	0.68	0.64	0.66	0.20	0.48	0.34
S-SVDD ψ_2	0.91	0.84	0.77	0.84	0.21	0.61	0.41	0.40	0.52	0.46
S-SVDD ψ_3	0.92	0.85	0.73	0.83	0.35	0.62	0.49	0.37	0.16	0.27
SVDD	0.94	0.91	0.84	0.89	0.31	0.80	0.55	0.53	0.66	0.59
E-SVDD	0.89	0.84	0.86	0.86	0.30	0.00	0.15	0.00	0.00	0.00
GE-SVDD	0.91	0.88	0.85	0.88	0.26	0.81	0.54	0.56	0.66	0.61
OC-SVM	0.45	0.65	0.66	0.59	0.27	0.63	0.45	0.51	0.58	0.54
GE-SVM	0.92	0.90	0.86	0.89	0.39	0.91	0.65	0.54	0.67	0.61
Non-support-vector-based methods										
SOM	0.91	0.84	0.88	0.88	0.06	0.87	0.47	0.46	0.32	0.39
Gaussian	0.97	0.86	0.80	0.88	0.33	0.50	0.42	0.47	0.59	0.53
GE-OC-ELM	0.99	0.89	0.78	0.88	0.48	0.81	0.65	0.38	0.55	0.47

covariance of the class. Regularization term ψ_1 , which uses all training samples to describe the covariance of the class, also performed well. Both of these regularization terms produced 2 out of 6 best results in the linear case. We also noticed that ES-SVDD without any regularization term performs the worst as compared to ES-SVDD with regularization terms.

For non-linear ES-SVDD, the regularization terms ψ_1 and Υ_3 resulted in the best results for most of the datasets. However, ψ_1 is also noticed to perform worse than the others in a few datasets. ψ_1 uses all target training samples in describing the covariance of the class without taking into account the concentration matrix. In Υ_3 , the λ values take the values of α_i values of the boundary samples and zero for non-boundary samples. In the non-linear case for high dimensional datasets, we notice that using all the training data for describing the covariance of the data in a projected space, with or without using the concentration matrix (i.e., ψ_1 or Υ_1), yielded the best results for ES-SVDD.

We further notice that by considering the class variance taking into account the concentration matrix in the

regularization term, ES-SVDD performed better in most datasets as compared to the regularization terms without considering the concentration matrix. Overall Υ_2 is found to be more robust than other regularization terms. Hence, we recommend to use samples belonging to the boundary and outside the boundary to describe the covariance of the class while taking into account the concentration matrix.

We also show the performance of the proposed ES-SVDD and the recently proposed S-SVDD on the test set after every training iteration for the linear and non-linear cases. We compare the performances of these methods with different regularization terms Υ and ψ . The average Gmean value is calculated for each iteration over the 5 test splits for the different datasets, see Figures 1-6.

It can clearly be seen from the figures that for both the linear and non-linear methods, ES-SVDD achieves its best performance much earlier than the recently proposed counterpart S-SVDD. This is not surprising, because the ellipsoidal description can fit a larger variety of data distributions,

Seeds (Kama)

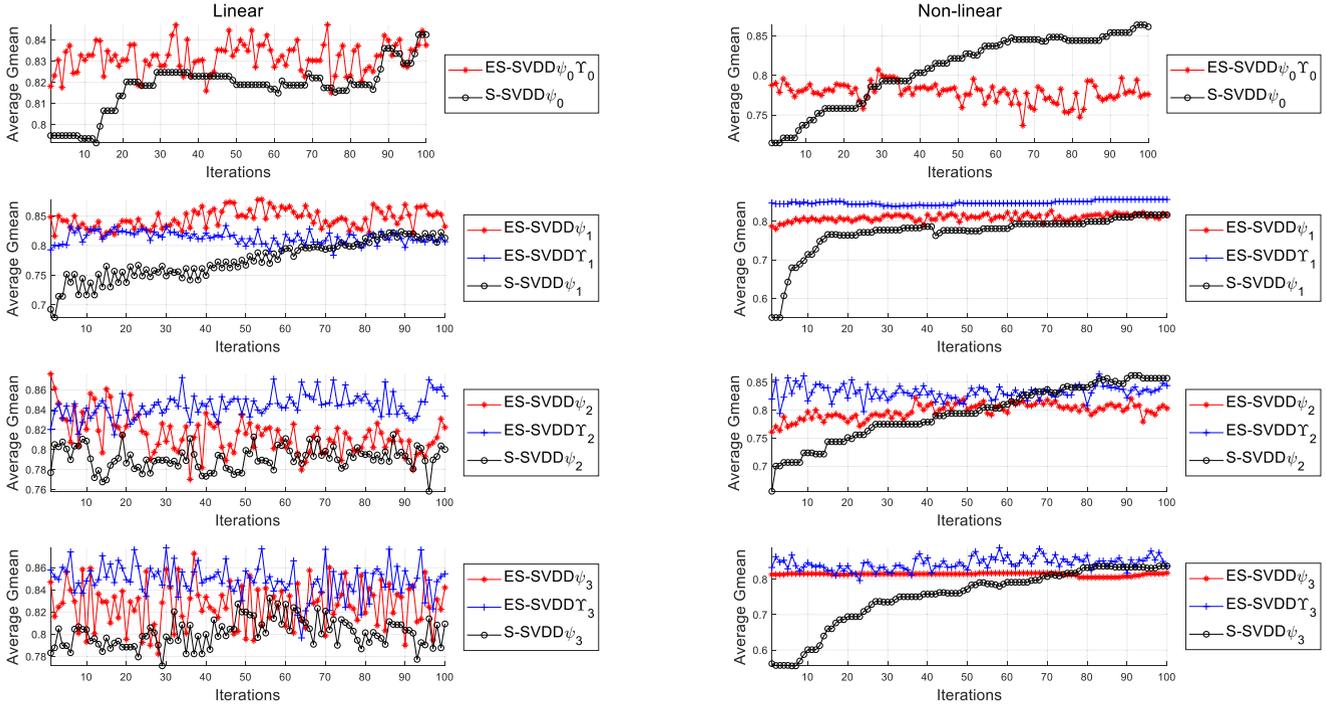


FIGURE 1. Comparison of different regularization terms for ES-SVDD and S-SVDD on dataset S-K.

Qualitative bankruptcy (bankruptcy)

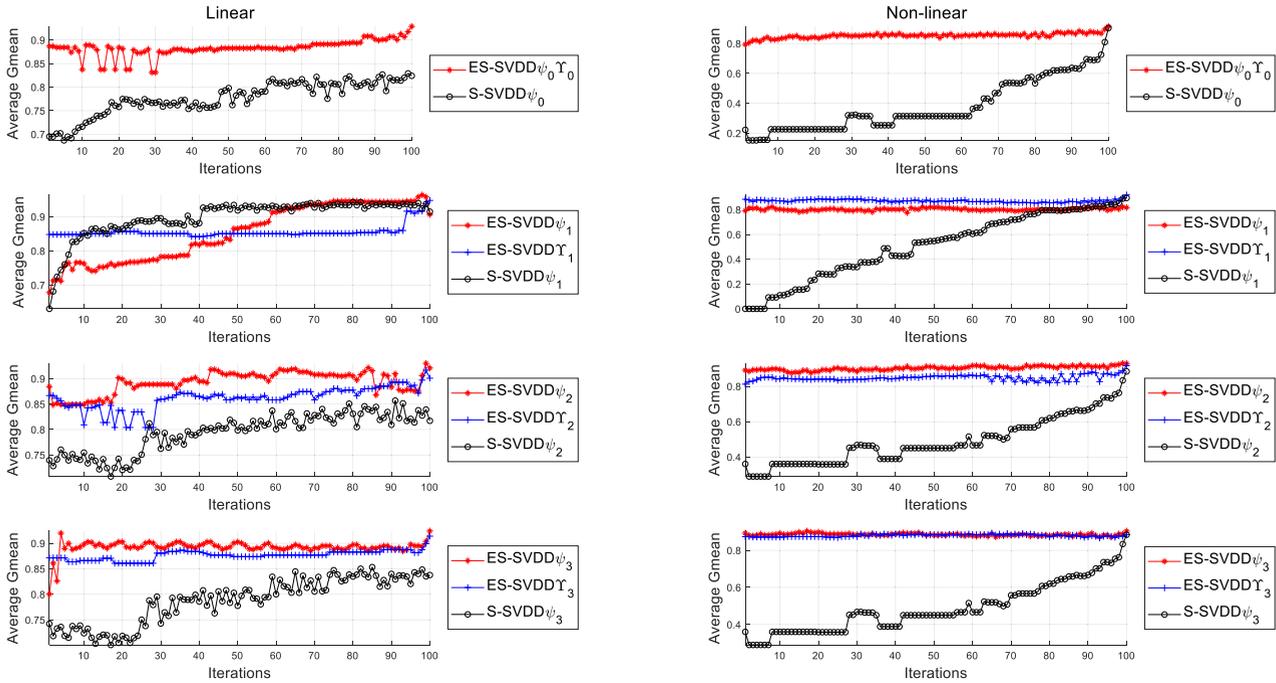


FIGURE 2. Comparison of different regularization terms for ES-SVDD and S-SVDD on dataset QB-B.

while the optimal spherical description gets successful only after the data variance for different dimensions has been equalized. Using the ellipsoidal data description in the

proposed method makes it converge faster to an optimal solution. We also notice that for high dimensional datasets ES-SVDD ψ_1 and ES-SVDD γ_1 are more stable as compared

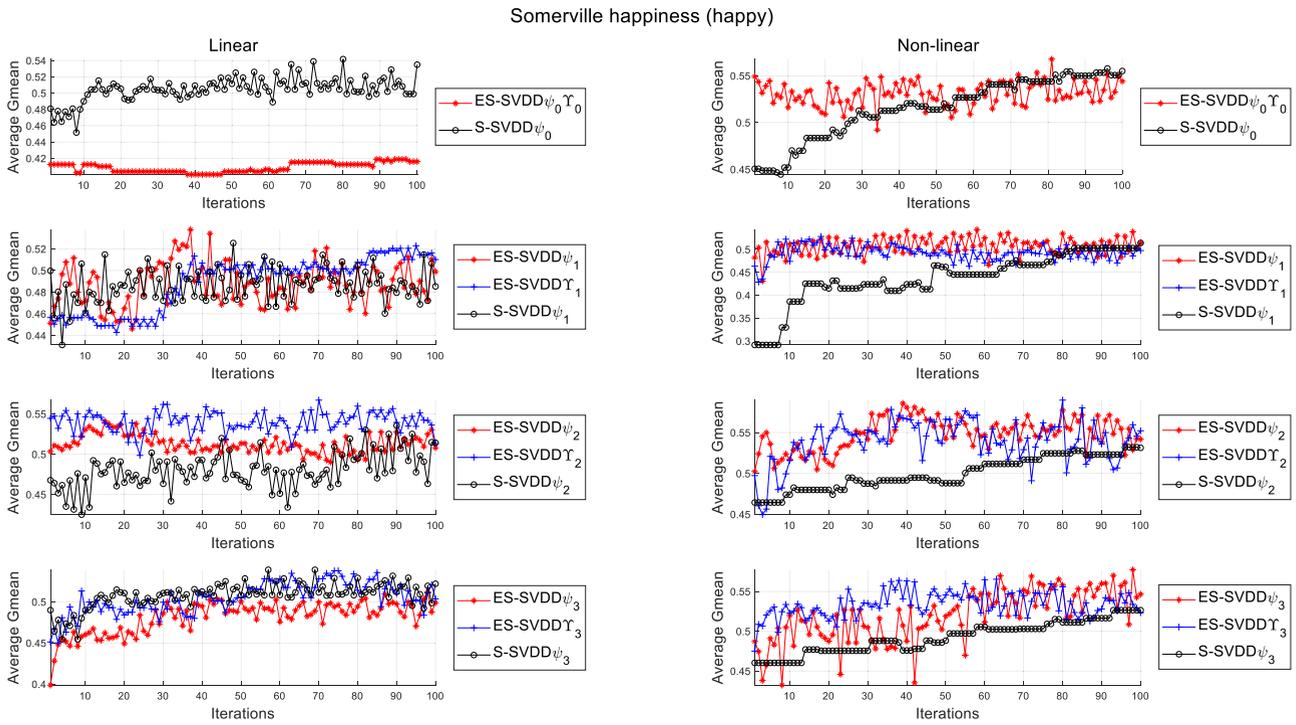


FIGURE 3. Comparison of different regularization terms for ES-SVDD and S-SVDD on dataset SH-H.

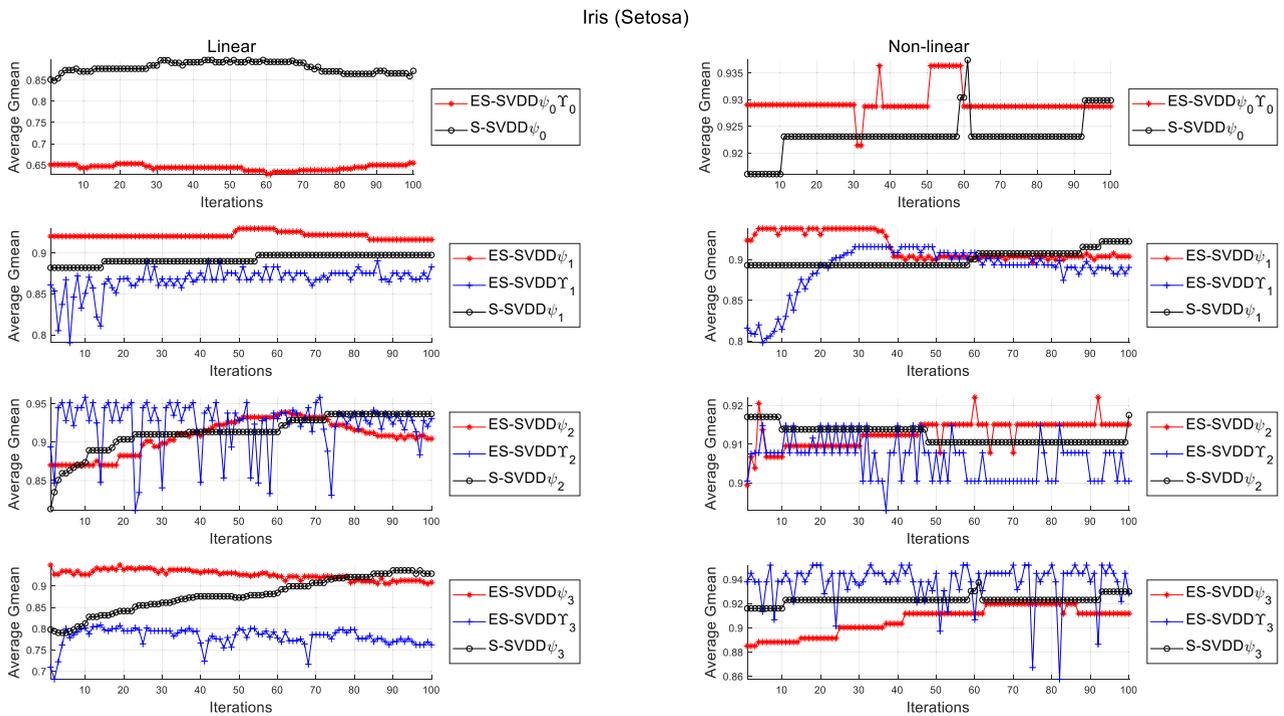


FIGURE 4. Comparison of different regularization terms for ES-SVDD and S-SVDD on dataset I-S.

to the other proposed linear and non-linear methods. Overall, the trend of faster convergence and higher stability in terms of producing consistent results for different range of

iterations for ES-SVDD can be observed both in the linear and non-linear methods for all regularization terms in the majority of the cases.

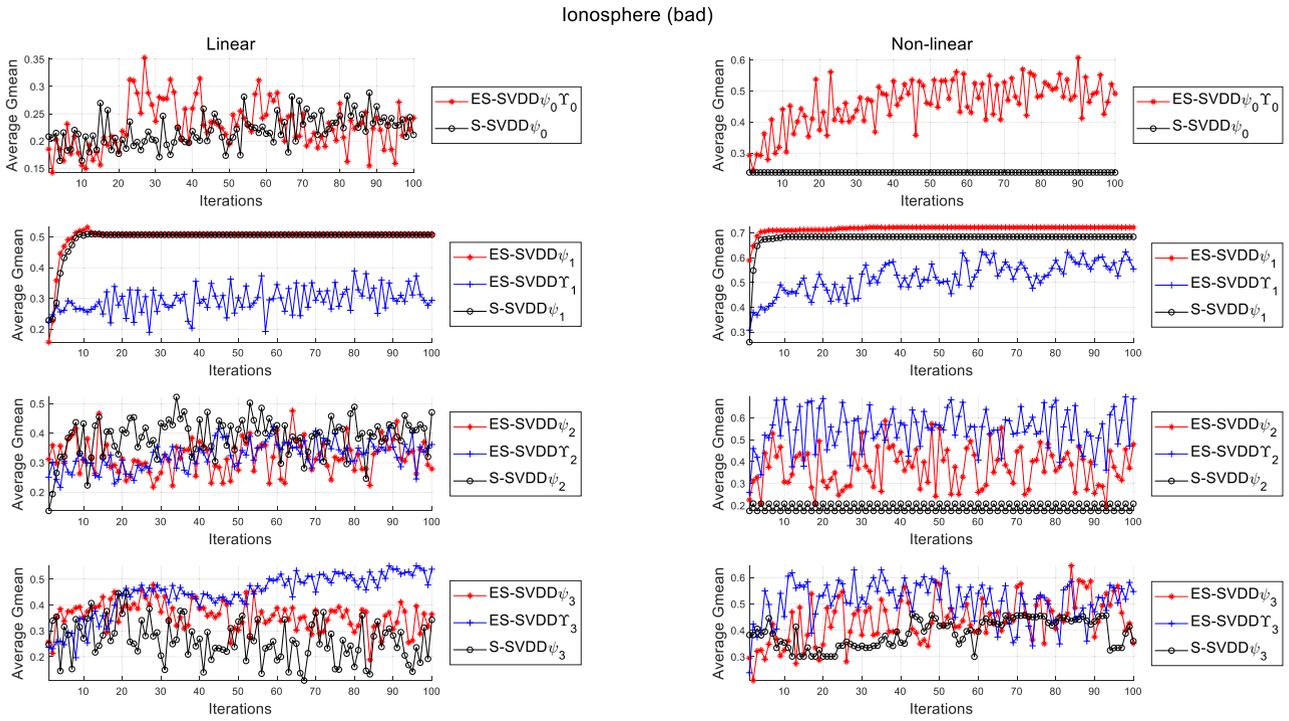


FIGURE 5. Comparison of different regularization terms for ES-SVDD and S-SVDD on dataset IS-B.

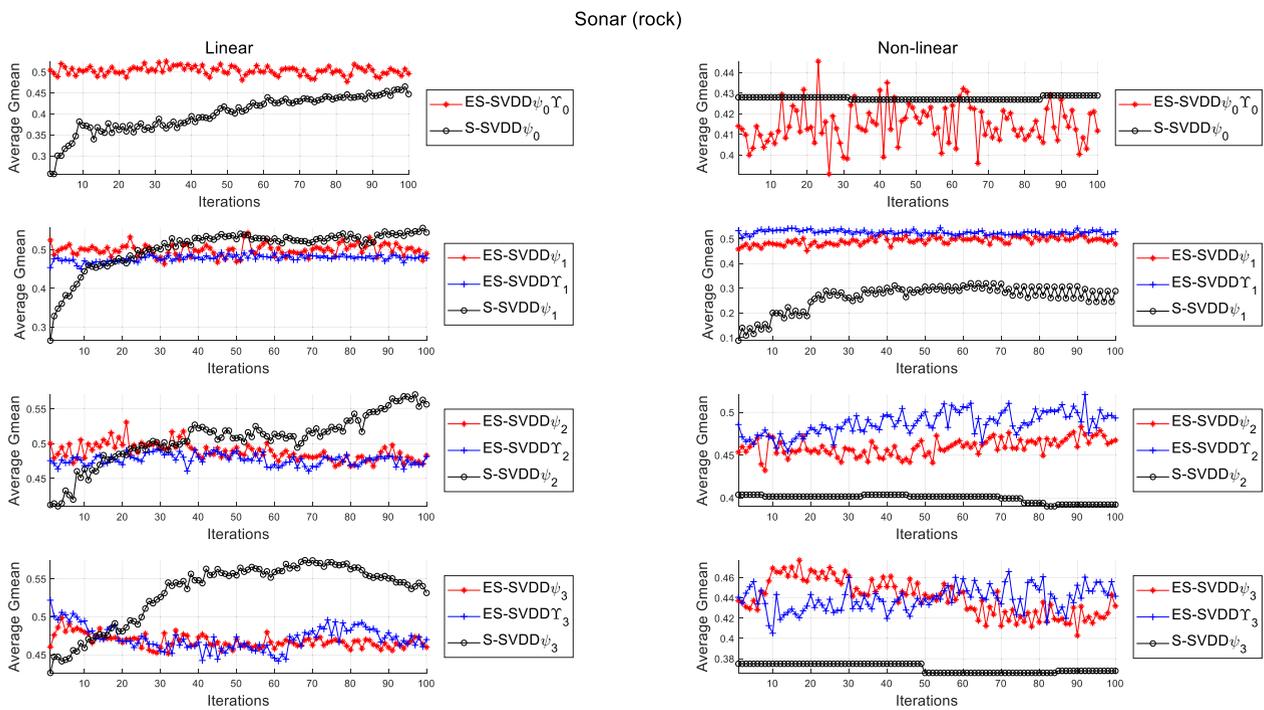


FIGURE 6. Comparison of different regularization terms for ES-SVDD and S-SVDD on dataset SR-R.

V. CONCLUSION

In this paper, a novel method, ES-SVDD, for one-class classification is proposed. The proposed method projects the data from an input feature space to a new optimized subspace

suitable for one-class classification. The proposed method generalizes S-SVDD for a hypersphere by using ellipsoidal data description. We proposed different regularization terms along with linear and non-linear formulations of the method.

In most cases, the proposed ES-SVDD variants outperform the competing SV-based methods and converge faster than in the case of data description without ellipsoidal encapsulation.

In the future, we intend to use other kernel types in the non-linear case of ES-SVDD. We also plan to devise a strategy for early exit in the training process to reduce the training time. We will also experiment with finetuning hyperparameters according to different criteria, such as area under receiver operating characteristic curve. Additionally, we plan to formulate and implement a neural network based version of the proposed method and compare its performance with deep neural networks.

REFERENCES

- [1] T. Hastie, R. Tibshirani, and J. Friedman, "Unsupervised learning," in *The Elements of Statistical Learning*. New York, NY, USA: Springer, 2009, pp. 485–585.
- [2] S. B. Kotsiantis, I. Zaharakis, and P. Pintelas, "Supervised machine learning: A review of classification techniques," *Emerg. Artif. Intell. Appl. Comput. Eng.*, vol. 160, no. 1, pp. 3–24, 2007.
- [3] K. Fukunaga, *Introduction to Statistical Pattern Recognition*. Amsterdam, The Netherlands: Elsevier, 2013.
- [4] Y. Yang, C. Hou, Y. Lang, G. Yue, and Y. He, "One-class classification using generative adversarial networks," *IEEE Access*, vol. 7, pp. 37970–37979, 2019.
- [5] A. Iosifidis, V. Mygdalis, A. Tefas, and I. Pitas, "One-class classification based on extreme learning and geometric class information," *Neural Process. Lett.*, vol. 45, no. 2, pp. 577–592, Apr. 2017.
- [6] V. Mygdalis, A. Iosifidis, A. Tefas, and I. Pitas, "Laplacian one class extreme learning machines for human action recognition," in *Proc. IEEE 18th Int. Workshop Multimedia Signal Process. (MMSP)*, Sep. 2016, pp. 1–5.
- [7] D. M. J. Tax, "One-class classification: Concept learning in the absence of counter-examples," *ASCI dissertation*, Delft Univ. Technol., Delft, The Netherlands, 2001, vol. 65, pp. 1–190.
- [8] K. Lee, D.-W. Kim, K. H. Lee, and D. Lee, "Density-induced support vector data description," *IEEE Trans. Neural Netw.*, vol. 18, no. 1, pp. 284–289, Jan. 2007.
- [9] C. Fraley and A. E. Raftery, "Model-based clustering, discriminant analysis, and density estimation," *J. Amer. Stat. Assoc.*, vol. 97, no. 458, pp. 611–631, Jun. 2002.
- [10] S. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Inf. Theory*, vol. 28, no. 2, pp. 129–137, Mar. 1982.
- [11] T. Kohonen and P. Somervuo, "Self-organizing maps of symbol strings," *Neurocomputing*, vol. 21, nos. 1–3, pp. 19–30, Nov. 1998.
- [12] B. Schölkopf, R. Williamson, A. Smola, and J. Shawe-Taylor, "SV estimation of a distribution's support," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 12, 1999, pp. 1–7.
- [13] D. M. J. Tax and R. P. W. Duin, "Support vector data description," *Mach. Learn.*, vol. 54, no. 1, pp. 45–66, Jan. 2004.
- [14] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft, "Deep one-class classification," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 4393–4402.
- [15] J. Zhao, H. Lin, Q. Chen, X. Huang, Z. Sun, and F. Zhou, "Identification of egg's freshness using NIR and support vector data description," *J. Food Eng.*, vol. 98, no. 4, pp. 408–414, Jun. 2010.
- [16] H. Lee and W. Chung, "Terrain classification for mobile robots on the basis of support vector data description," *Int. J. Precis. Eng. Manuf.*, vol. 19, no. 9, pp. 1305–1315, Sep. 2018.
- [17] V. Mygdalis, A. Iosifidis, A. Tefas, and I. Pitas, "Graph embedded one-class classifiers for media data classification," *Pattern Recognit.*, vol. 60, pp. 585–595, Dec. 2016.
- [18] J. Wang, W. Liu, K. Qiu, H. Xiong, and L. Zhao, "Dynamic hypersphere SVDD without describing boundary for one-class classification," *Neural Comput. Appl.*, vol. 31, no. 8, pp. 3295–3305, Aug. 2019.
- [19] T. Kenaza, K. Bennaceur, and A. Labed, "An efficient hybrid SVDD/clustering approach for anomaly-based intrusion detection," in *Proc. 33rd Annu. ACM Symp. Appl. Comput. (SAC)*, 2018, pp. 435–443.
- [20] V. Mygdalis, A. Iosifidis, A. Tefas, and I. Pitas, "Semi-supervised subclass support vector data description for image and video classification," *Neurocomputing*, vol. 278, pp. 51–61, Feb. 2018.
- [21] M. Rahmanianesh, J. A. Nasiri, S. Jalili, and N. M. Charkari, "Adaptive three-phase support vector data description," *Pattern Anal. Appl.*, vol. 22, no. 2, pp. 491–504, May 2019.
- [22] B. Liu, Y. Xiao, L. Cao, Z. Hao, and F. Deng, "SVDD-based outlier detection on uncertain data," *Knowl. Inf. Syst.*, vol. 34, no. 3, pp. 597–618, Mar. 2013.
- [23] A. Banerjee, P. Burlina, and R. Meth, "Fast hyperspectral anomaly detection via SVDD," in *Proc. IEEE Int. Conf. Image Process.*, vol. 4, 2007, p. IV-101.
- [24] Y.-H. Liu, Y.-C. Liu, and Y.-J. Chen, "Fast support vector data descriptions for novelty detection," *IEEE Trans. Neural Netw.*, vol. 21, no. 8, pp. 1296–1313, Aug. 2010.
- [25] F. Camci and R. B. Chinnam, "General support vector representation machine for one-class classification of non-stationary classes," *Pattern Recognit.*, vol. 41, no. 10, pp. 3021–3034, Oct. 2008.
- [26] Y. Forghani, S. Effati, H. S. Yazdi, and R. S. Tabrizi, "Support vector data description by using hyper-ellipse instead of hyper-sphere," in *Proc. 1st Int. eConf. Comput. Knowl. Eng. (ICCKE)*, Oct. 2011, pp. 22–27.
- [27] M. GhasemiGol, R. Monsefi, and H. S. Yazdi, "Ellipse support vector data description," in *Proc. Int. Conf. Eng. Appl. Neural Netw.* Berlin, Germany: Springer, 2009, pp. 257–268.
- [28] K. Wang and H. Xiao, "Ellipsoidal data description," *Neurocomputing*, vol. 238, pp. 328–339, May 2017.
- [29] M. Ghasemigol, R. Monsefi, and H. Sadoghi-Yazdi, "Intrusion detection by ellipsoid boundary," *J. Netw. Syst. Manage.*, vol. 18, no. 3, pp. 265–282, Sep. 2010.
- [30] E. Rimon and S. P. Boyd, "Obstacle collision detection using best ellipsoid fit," *J. Intell. Robot. Syst.*, vol. 18, no. 2, pp. 105–126, 1997.
- [31] K. Wang, H. Xiao, and Y. Fu, "Ellipsoidal support vector data description in kernel PCA subspace," in *Proc. 3rd Int. Conf. Digit. Inf. Process., Data Mining, Wireless Commun. (DIPDMWC)*, Jul. 2016, pp. 13–18.
- [32] F. Sohrab, J. Raitoharju, M. Gabbouj, and A. Iosifidis, "Subspace support vector data description," in *Proc. 24th Int. Conf. Pattern Recognit. (ICPR)*, Aug. 2018, pp. 722–727.
- [33] N. Kwak, "Nonlinear projection trick in kernel methods: An alternative to the kernel trick," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 24, no. 12, pp. 2113–2119, Dec. 2013.
- [34] R. Chen, N. Sun, X. Chen, M. Yang, and Q. Wu, "Supervised feature selection with a stratified feature weighting method," *IEEE Access*, vol. 6, pp. 15087–15098, 2018.
- [35] H. Luo and J. Han, "Trace ratio criterion based large margin subspace learning for feature selection," *IEEE Access*, vol. 7, pp. 6461–6472, 2019.
- [36] M. U. Chaudhry and J.-H. Lee, "Feature selection for high dimensional data using Monte Carlo tree search," *IEEE Access*, vol. 6, pp. 76036–76048, 2018.
- [37] Z. Zhao, L. Wang, H. Liu, and J. Ye, "On similarity preserving feature selection," *IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 3, pp. 619–632, Mar. 2013.
- [38] N. Zhou, Y. Xu, H. Cheng, J. Fang, and W. Pedrycz, "Global and local structure preserving sparse subspace learning: An iterative approach to unsupervised feature selection," *Pattern Recognit.*, vol. 53, pp. 87–101, May 2016.
- [39] R. Vijayanand and D. Devaraj, "A novel feature selection method using whale optimization algorithm and genetic operators for intrusion detection system in wireless mesh network," *IEEE Access*, vol. 8, pp. 56847–56854, 2020.
- [40] Q. Gu, Z. Li, and J. Han, "Joint feature selection and subspace learning," in *Proc. 22nd Int. Joint Conf. Artif. Intell.*, 2011, pp. 1294–1299.
- [41] M. Bacher, I. Ben-Gal, and E. Shmueli, "Subspace selection for anomaly detection: An information theory approach," in *Proc. IEEE Int. Conf. Sci. Electr. Eng. (ICSEE)*, Nov. 2016, pp. 1–5.
- [42] H. V. Nguyen, E. Muller, and K. Bohm, "4S: Scalable subspace search scheme overcoming traditional apriori processing," in *Proc. IEEE Int. Conf. Big Data*, Oct. 2013, pp. 359–367.
- [43] K. B. Petersen and M. S. Pedersen. (Nov. 2012). *The Matrix Cookbook, Version 20121115*. [Online]. Available: <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
- [44] D. Dua and C. Graff. (2017). *UCI Machine Learning Repository*. [Online]. Available: <http://archive.ics.uci.edu/ml>

- [45] J. Laaksonen and T. Honkela, *Advances in Self-Organizing Maps: 8th International Workshop, WSOM 2011, Espoo, Finland, June 13-15, 2011. Proceedings*, vol. 6731. Berlin, Germany: Springer, 2011.
- [46] D. M. J. Tax and P. Juszczak, "Kernel whitening for one-class classification," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 17, no. 3, pp. 333–347, May 2003.
- [47] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," *ACM Trans. Intell. Syst. Technol.*, vol. 2, no. 3, pp. 27:1–27:27, 2011. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [48] D. M. J. Tax, "A MATLAB toolbox for data description, outlier and novelty detection version 1.7.5," Delft Univ. Technol., Delft, The Netherlands, Tech. Rep., Apr. 2010.



FAHAD SOHRAB (Graduate Student Member, IEEE) received the B.S. degree (*cum laude*) in telecommunication engineering from the National University of Computer and Emerging Sciences, Peshawar Pakistan, in August 2012, and the master's degree from Sabanci University, Istanbul, Turkey. He is currently pursuing the Ph.D. degree with Tampere University, Finland. In the fall of 2013, he joined the Computer Vision and Pattern Analysis Laboratory, Sabanci University. He was

also affiliated with the Pattern Recognition Laboratory, Delft University of Technology, The Netherlands, during the second year of his Master's studies. He is currently with the Signal Analysis and Machine Intelligence Group, Tampere University. His research interests include machine learning, pattern recognition, subspace learning, and one-class classification.



JENNI RAITOHARJU (Member, IEEE) received the Ph.D. degree from the Tampere University of Technology, Finland, in 2017. Since then, she has worked as a Postdoctoral Research Fellow at the Faculty of Information Technology and Communication Sciences, Tampere University, Finland. In 2019, she started working as a Senior Research Scientist at the Finnish Environment Institute, Jyväskylä, Finland, after receiving Academy of Finland Postdoctoral Researcher funding for 2019–2022. She has coauthored 15 journal articles and 27 papers in international conferences. Her research interests include machine learning and pattern recognition methods along with applications in biomonitoring and autonomous systems. She has been the Chair of the Young Academy Finland, since 2019.



ALEXANDROS IOSIFIDIS (Senior Member, IEEE) received the B.Sc. degree in electrical and computer engineering and the M.Sc. degree with a specialisation in mechatronics from the Democritus University of Thrace, Greece, in 2008 and 2010, respectively, and the Ph.D. degree in computer science from the Aristotle University of Thessaloniki, Greece, in 2014. He is currently an Associate Professor of machine learning with the Department of Engineering, Aarhus University,

Denmark. Before, he joined Aarhus University, he held a Postdoctoral Researcher positions at the Aristotle University of Thessaloniki and at the Tampere University of Technology, Finland, where he was an Academy of Finland Postdoctoral Research Fellow. He has contributed in more than 20 Research and Development projects financed by EU, Finnish and Danish funding agencies and companies. He has (co)authored 65 articles in international journals and 85 papers in international conferences proposing novel Machine Learning techniques and their application in a variety of problems. He has served as an Officer of the Finnish IEEE Signal Processing—Circuits and Systems Chapter from 2016 to 2018. His research interests include neural networks and statistical machine learning finding applications in computer vision, financial engineering, and graph mining problems. He is also a member of the EURASIP Technical Area Committee on Visual Information Processing. He serves as an Area/Associate Editor in *Neurocomputing*, *Signal Processing: Image Communications*, *IEEE ACCESS*, and *BMC Bioinformatics* journals.



MONCEF GABBOUJ (Fellow, IEEE) received the B.S. degree in electrical engineering from Oklahoma State University, Stillwater, OK, USA, in 1985, and the M.S. and Ph.D. degrees in electrical engineering from Purdue University, West Lafayette, IN, USA, in 1986 and 1989, respectively. He was an Academy of Finland Professor, from 2011 to 2015. He is currently a Professor of signal processing with the Faculty of Information Technology and Communication Sciences, Tampere University, Finland. He guided 40 Ph.D. students and has published 650 articles. His current research interests include multimedia content-based analysis, indexing and retrieval, machine learning, nonlinear signal and image processing and analysis, voice conversion, and video processing and coding. He is a member of the Academia Europaea and the Finnish Academy of Science and Letters. He organized several tutorials and special sessions for major IEEE conferences and EUSIPCO. He is the past Chairman of the IEEE CAS TC on DSP and a Committee Member of the IEEE Fourier Award for Signal Processing. He has served as an Associate Editor and a Guest Editor for many IEEE international journals and a Distinguished Lecturer for the IEEE CASS.